# Chapter 16
# Identifying the Reaction Mechanisms of Inteins with QM/MM Multiscale Methods

1
2
3

**Philip T. Shemella and Saroj K. Nayak**

4

**Abstract** With a series of quantum mechanical calculations ranging from gas phase, to an implicit solvent scheme, to combined quantum/classical simulations, we have provided insight into some of the key steps of intein reactions. These studies may be exploited for many applications involving inteins including molecular switches and sensors as well as controlled drug delivery.

5
6
7
8
9

AQ1

## 1 Introduction

10

### *1.1 Computational Background*

11

Nestled between experiment and pure theory, computational chemistry has become an integral tool for researchers working in physics, chemistry, and biology, as well as nanotechnology and biotechnology. Computer simulations allow the researcher to access states both visible and invisible to experiment, and make predictions based on this knowledge. A chemical reaction may be quantified by the amount of reactants, the amount of products, and the time elapsed. To explain a mechanism and molecular structure and energies on the atomic level, computational methods are important.

12
13
14
15
16
17
18
19

The field of computational chemistry spans many length and time scales. To simulate protein folding, which requires an extremely long simulation trajectory, amino acids may be "coarse-grained," where the atomic description of each side chain is aggregated into a composite value. To achieve long trajectories this approximation as well as others are essential. However, to calculate the $pK_a$ of a side chain or the chemical shifts *via* nuclear magnetic resonance (NMR), not only will an atomic level description be necessary, but also a method that can calculate observable properties from first principles is often required.

20
21
22
23
24
25
26
27

AQ2

P.T. Shemella and S.K. Nayak
Rensselaer Polytechnic Institute, 110 8th St., Troy, NY 12180, USA
e-mail: nayaks@rpi.edu

439

The energies associated with bond breakage and formation are an essential 28
property for an enzymatic processes. For example, a change in energy barrier of 29
~1.4 kcal/mol at room temperature corresponds to an order of magnitude change 30
in the reaction rate. States observed at equilibrium may be predicted based on 31
relative energies between structures. To computationally access the energy of the 32
system, and to do so not only for equilibrium structures but also for transition states, 33
first principles electronic structure calculations are required. Using an all-electron 34
method, the electron orbitals are considered variable and flexible, and they depend 35
on neighboring atoms and environment. This is important because the chemistry 36
at transition states may vary greatly from equilibrium structures: instead of four 37
bonds, carbon atoms may have three or five bonds during a chemical reaction. 38
Transition states are where quantum mechanical principles dominate. By solving 39
the Schrödinger equation for all electrons, and relaxing their orbital positions and 40
therefore allowing the electron density to vary, an accurate description of the system 41
can be obtained that is useful for understanding fundamental chemistry both near 42
and far from equilibrium. 43

## 1.2   Inteins Background                                                        44

Protein splicing involves the autocatalytic release of a peptide segment, termed an 45
intein, with the joining of two flanking protein sequences (exteins) [1, 2]. Inteins 46
are autocatalytic proteins that exist in all three domains of life. Experiments have 47
identified key reaction steps in protein splicing whereas sequence comparisons have 48
revealed the conserved amino acids required for this reaction. Figure 1 shows a 49
schematic for conserved intein residues and their corresponding block (C or N) des- 50
ignation. Experimental mutational studies have been carried out to further control 51
the protein splicing reaction [3, 4]. For example, by mutating the first residue at the 52
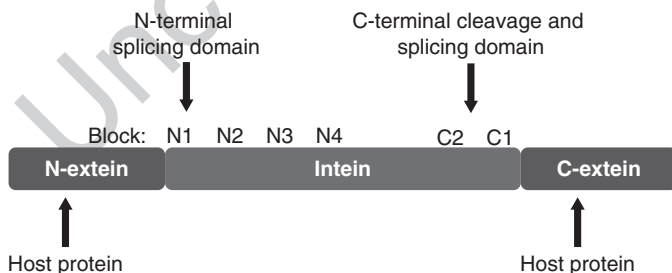N-terminus (N1 block) of the intein from Cys to Ala (N1-Cys1Ala), the first step of 53



**Fig. 1** Schematic intein and N- and C-exteins. Splicing motifs contain highly conserved amino acids, such as N1-Cys1, N3-His10, C2-Asp5, and C1-His7, C1-Asn8, C1-Cys+1

the splicing reaction, namely the N-terminal N–S shift,[1] is inhibited, thus isolating 54
the C-terminal cleavage reaction [5]. Mutation schemes that control the reaction 55
rate and/or the specific products could be exploited in many biotechnological appli- 56
cations such as bioseparations [6, 7], drug development [8], and molecular sensors 57
[9, 10]. 58

## 2    Methods 59

### 2.1    Computational Methodology 60

In order to obtain an atomic-level understanding on the reaction mechanisms as 61
well as on the effect of mutation on the reaction barrier, we have carried out de- 62
tailed quantum mechanical simulations on intein C-terminal cleavage reactions. We 63
describe pH dependent C-terminal cleavage calculations for the *Mtu* recA intein; 64
performed with semi-empirical, QM gas phase, QM implicit solvent, and combined 65
QM/MM calculations [11–13]. Harnessing the C-terminal cleavage reaction may al- 66
low for an intein-based delivery device, where the reaction is triggered by a certain 67
stimulus. 68

Our computational results indicate that certain mutations either inhibit or en- 69
hance specific reaction steps of the overall splicing reaction, a conclusion that is 70
consistent with experiment. With quantum mechanical simulations, intermediate 71
states may be isolated and studied in the context of altering the molecular triggers 72
and inhibitors that impact protein splicing with inteins. The ability to study precur- 73
sor, intermediate, and post-reaction product states is extremely useful and carried 74
out with first principles methods. 75

### 2.2    Quantum Mechanical (QM) Methods 76

First principles density functional theory (DFT) [14, 15] was used to study intein 77
C-terminal cleavage; in particular, the Becke three-parameter hybrid functional, 78
B3LYP [16]. This hybrid method combines exchange terms from the Local Spin 79
Density Approximation (LSDA), Hartree–Fock (HF), and Becke's (B88) exchange 80
[17] with the correlation functionals from Lee, Yang, and Parr (LYP) [18] as well as 81
that from the LSDA [19]. The exchange (X) and correlation (C) energy is written as 82
$E_{XC}^{B3LYP}$, where 83

$$E_{XC}^{B3LYP} = (1-a)E_X^{LSDA} + aE_X^{HF} + b\Delta E_X^{Becke} + E_C^{LSDA} + c\Delta E_C^{LYP}, \quad (1)$$

---

[1] Atoms are annotated with one letter, i.e., H = hydrogen. Amino acids are annotated with three letters, i.e., His = histidine.

and the coefficients were optimized to match extensive molecular data sets ($a$ = 84
0.20, $b$ = 0.72, and $c$ = 0.81) [16]. Implemented with Gaussian code [20], this 85
hybrid gradient-corrected method is considered one of the most accurate exchange- 86
correlation functionals and has been used with great success in other biological 87
systems [21, 22]. Calculations with post-Hartree–Fock Møller–Plesset perturbation 88
theory (MP2) [23–26] were conducted to test the accuracy of the B3LYP method 89
for this system, and the energy barrier calculations were consistent [12]. 90

The first term in the hybrid method is $E_x^{LDA}$, which is the local density approxi- 91
mation (LDA) exchange term. $E_x^{HF}$ is Hartree–Fock exchange integral, which is an 92
exact quantity for electron spin exchange. Becke's B88 exchange term [17] is based 93
on empirical results, and is written as, 94

$$E_x^{Becke}[\rho(\mathbf{r})] = -\beta \int d\mathbf{r}\rho(\mathbf{r})^{4/3} \frac{\alpha^2}{(1 + 6\beta \sinh^{-1} \alpha)} \qquad (2)$$

where

$$\alpha = \frac{\mid \nabla\rho(\mathbf{r}) \mid}{\rho(\mathbf{r})^{4/3}}.$$

Found by matching molecular data sets, $\beta$ was found to be 0.0042 Hartree. Cor- 95
relation functionals are from the LDA [19] and from Lee, Yang, and Parr (LYP) 96
[18, 27], the latter based on an empirically determined model of the correlation en- 97
ergy of electrons in a helium atom. 98

Implemented with Gaussian code [20], this hybrid gradient-corrected method is 99
considered one of the most accurate exchange-correlation functionals and has been 100
used with great success in other biological systems [21, 22]. 101

We have used the double-$\zeta$ basis set, 6-31G(d,p), for geometry optimizations dur- 102
ing initial reaction path sampling [28], where the '6' represents six GTOs for core 103
electrons and the '31' represents split GTOs for valence electrons: specifically three 104
and one. Split-valence basis sets allow for a more accurate description of chemical 105
bonding due to increased flexibility to fit valence electrons into molecular orbitals, 106
and are the norm when using a Gaussian-type basis set. The '(d,p)' indicates that 107
we are using polarization functions that allow for a shift in the wave function away 108
from the atomic center. We have also used the triple-$\zeta$ basis set, 6-311++g(d,p), for 109
calculations of the local minima and transition states found with the first basis set 110
[29]. Diffuse functions for long range interactions are represented with a '+', and 111
are especially important for anions. Basis sets of similar size are typically used for 112
systems with similar number of electrons, and our test calculations as well as the 113
work of others have shown these basis sets to be sufficient for similar atom types 114
[21, 22]. 115

### 2.2.1 Implicit Solvent                                                                                                        116

One method for approximating the environmental electrostatic effect is to use an 117
implicit solvent. In this scheme, the active site is polarized by the dielectric medium 118

which is itself polarizable. The Polarizable Continuum Model (PCM) [30] was   119
used to simulate solvent effects in the detailed calculations. The numerical Integral   120
Equation Formalism [31] (IEFPCM) was used because it allows for interlocking   121
atomic spheres to represent the extent of the system in solution, which is important   122
for protons that are in between atoms during a chemical reaction and at or around   123
the energy barrier.   124

Non-dimensional dielectric constants are defined by $\varepsilon_r = \varepsilon_s/\varepsilon_0$, where $\varepsilon_0$ is the   125
vacuum permittivity and $\varepsilon_s$ is the static dielectric constant for the dielectric. For the   126
gas phase, $\varepsilon_r = 1$. For water, $\varepsilon_r = 78.39$. Geometry optimizations were performed   127
in implicit solvent and results are compared with gas phase calculations.   128

## 2.3   Classical Methods   129

Starting with the intein crystal structure for the *Mtu* RecA intein, ($\Delta\Delta$Ihh-CM, PDB   130
code 2IN8) [32], a product protein without exteins, N- and C-terminal exteins were   131
computationally added and then equilibrated with classical molecular dynamics   132
(MD) simulations. The N-extein sequence consisted of Ace-Val-Val-Lys-Asn-Lys   133
and the C-extein sequence consisted of Cys-Ser-Pro-Pro-Phe-Nme, both based on   134
the native extein sequences [33]. Ace and Nme were capping residues for the N and   135
C-terminal exteins, respectively. AMBER force field parameters [34] were imple-   136
mented with GROMACS code [35]. MD simulations were carried out for 4 ns (0.5   137
ns equilibration, 3.5 ns production run) with temperature T = 298 K, pressure =   138
1bar, and number of water molecules = 9548 for Cys and 9549 for Met systems.   139

## 2.4   Multiscale (QM/MM) Methods   140

The QM/MM layering method involves treating the protein active site and criti-   141
cal solvent molecules with first principles methods while treating the remaining   142
full-protein system with classical force fields [36]. The classical periodic system   143
was trimmed down to include the protein (intein and exteins) as well as all interior   144
waters and those exterior water molecules within a range of 7.0 Å to the protein   145
surface (as a reference, the lone protein is roughly shaped like an oblate spheroid   146
and approximately $25 \times 35 \times 35$ Å$^3$). All atoms were relaxed, and each calculation   147
included at least 6,500 atoms. The full-protein plus solvent system, termed the real   148
system, was treated only with the MM method. Within the real system, the active   149
site model system was partitioned, and was treated independently by QM and MM   150
methods. Dangling bonds that were introduced by partitioning the model system   151
were then passivated with hydrogen atoms. With normal QM/MM energy calcula-   152
tions and geometry optimizations, protein and solution outside the model system   153
was only included as a mechanical perturbation. For this reason, it is critical that   154

the model system should include protein segments and solution molecules that are    155
interacting electrostatically. The combined Hamiltonian may be written:              156

$$E^{QM:MM}_{ONIOM} = E^{QM}_{Model} - E^{MM}_{Model} + E^{MM}_{Real} \tag{3}$$

For the smaller model system, $E^{QM}_{Model}$ is the energy calculated with quantum me-    157
chanical methods while the energy calculated by classical molecular mechanics        158
methods is given by $E^{MM}_{Model}$. The real system (full protein + solvent) energy is cal-    159
culated with MM methods and is given by $E^{MM}_{Real}$. In addition to the mechanical    160
perturbation on the QM Hamiltonian, the electrostatic contribution from the partial  161
charges of the MM region can be included as a perturbation on the QM Hamiltonian     162
[37]. The QM/MM formalism has been used with success in previous work [38–40].       163
Typically, we report $E^{QM}_{Model}$, which represents the QM active site energy. The other    164
energy terms, including the combined $E^{QM:MM}_{ONIOM}$ involves classical parameters de-    165
termined for equilibrium structures that have no relevance to the energies of bond   166
forming and breaking at transition states.                                           167

### 2.4.1  Charge Embedding                                                          168

In addition to the mechanical perturbation on the QM Hamiltonian, the electrostatic  169
contribution from the partial charges of the MM region can be included as a pertur-  170
bation on the QM Hamiltonian. For this scheme the partial charges are those used     171
in the MM calculation and are scaled by the default manner where atoms bonded        172
to the inner-most four layers and atoms outside that threshold are not included [37]. 173
Typically, we report $E^{QM}_{Model}$, which represents the QM active site energy. The other    174
energy terms, including the combined $E^{QM:MM}_{ONIOM}$ involves classical parameters that    175
have no relevance to the energies of bond forming and breaking at transition states.  176

## 2.5  Geometry Minimization                                                        177

Due to the complexity of biomolecular reactions, a rigorous multidimensional         178
search over local conformational space is essentially required although not computa- 179
tionally feasible for large systems [41]. Due to the time expense for each calculation, 180
we have used the constant minimization procedure. For intermediate states along the  181
reaction path, one coordinate is constrained while the remaining system is relaxed.  182
The constrained internal coordinate, called the Asn cyclization distance, was the    183
atomic distance between the Asn side chain N atom and the carbonyl C of Asn on       184
the scissile peptide bond. In calculations with a hydronium ion ($H_3O^+$), the three    185
O–H bond distances were often constrained to 0.98 Å to avoid spontaneous proton      186
donation observed otherwise.                                                         187

# 3   Results

## 3.1   Non-essential Mutation

Once splicing was inhibited, the downstream Cys residue (which was the first amino acid of the C-terminal extein or C-extein) was found to be functionally unnecessary for the C-terminal cleavage mechanism. Interestingly, Wood et al. observed that this amino acid regulated the reaction rate but did not alter the mechanism [42]. Furthermore, since the CM was found to be exceedingly reactive at low pH values, Wood et al. [42] utilized Met, which was the native N-terminus of the protein that formed the C-extein sequence, to decrease the reaction rate by an order of magnitude. In this experiment, three proteins of various sizes were contrasted with only the Cys/Met C-extein mutation: Thymidylate synthase (31.5 kDa), Hfq Protein (18 kDa), and rh aFGF (14 kDa). For these proteins, the Cys to Met mutation resulted in a decrease of the reaction rate by a factor of 12.0, 5.0, and 7.8, respectively [42, 43]. Figure 2 shows a schematic of the intein precursor and products based on these results [10, 44], although the exact mechanisms that govern the splicing and cleavage reactions are not understood at the atomic level. In particular, the effect of the single amino acid mutation at C + 1, flanking the conserved C1: His7-Asn8 dipeptide at the intein terminus, on the reaction rate is not understood.
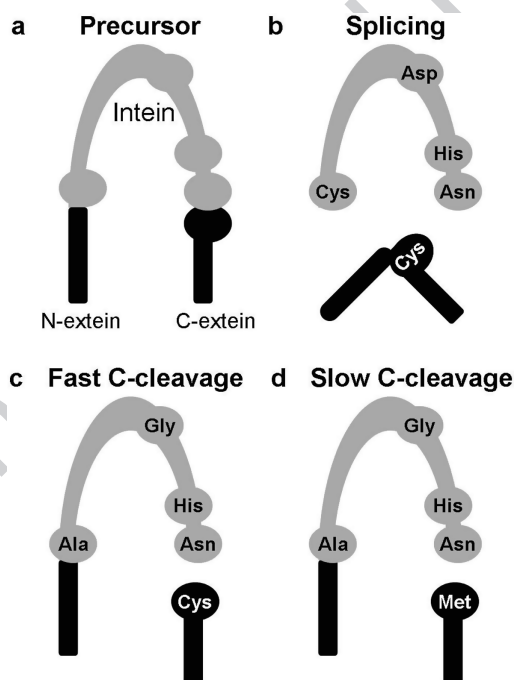


**Fig. 2**  Intein and extein precursor (**a**) and three possible reactions based on mutagenesis results: splicing product (**b**), and fast (**c**) and slow (**d**) C-terminal cleavage product

In order to obtain an atomic-level understanding of the effect of mutation on the    206
reaction barrier, detailed quantum mechanical calculations on the intein C-terminal   207
cleavage reaction have been carried out [12]. Simulations were based on both full     208
quantum mechanical molecular analysis as well as a hybrid quantum mechanics and      209
molecular mechanics (QM/MM) approach where the entire protein and solvent are        210
treated classically with parameterized force fields in a molecular mechanics (MM)    211
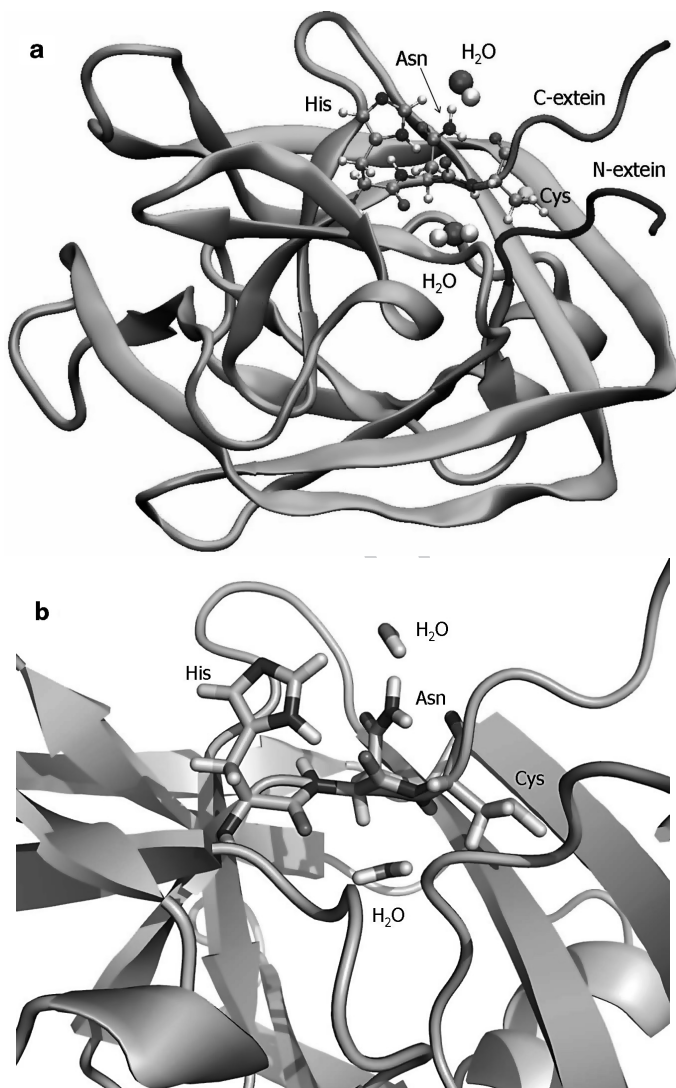calculation as shown in Fig. 3a. The 53 atom C-terminal catalytic site (C1-block:    212



**Fig. 3** The intein cleavage mutant (CM) crystal structure (PDB code 2IN8) with computation-ally added exteins (**a**). The C-terminal catalytic site (His–Asn–Cys + two water molecules) is highlighted (**b**)

His–Asn–Cys, or His–Asn–Xxx, where Xxx is an alternate amino acid) was treated    213
with quantum mechanics (QM) and is shown in Fig. 3b.    214

The computational energy barrier was smaller for the C-terminal sequence His–    215
Asn–Cys than for that of the His–Asn–Met mutant, consistent with experimental    216
observations [42, 43]. The difference in energy barrier between Cys/Met residues    217
was due to the difference in electron affinity of the amino acids. In addition to Cys    218
and Met, several other amino acids at the first C-extein position (C+1) were studied    219
here. The energy barrier for C-terminal cleavage, calculated with a larger model    220
system, is confirmed to match with that of the experiment.    221

## 3.2    Classical Protein System    222

Starting with the intein crystal structure for the *Mtu* recA intein, ($\Delta\Delta$Ihh-CM, PDB    223
code 2IN8) [32], a product protein without exteins, N- and C-terminal exteins were    224
computationally added and then equilibrated with classical molecular dynamics    225
(MD) simulations. The N-extein sequence consisted of Ace-Val-Val-Lys-Asn-Lys    226
and the C-extein sequence consisted of Cys-Ser-Pro-Pro-Phe-Nme, both based on    227
the native extein sequences [33]. Ace and Nme were capping residues for the N and    228
C-terminal exteins, respectively. AMBER force field parameters [34] were imple-    229
mented with GROMACS code [35]. MD simulations were carried out for 4 ns (0.5    230
ns equilibration, 3.5 ns production run) with temperature T = 298K, pressure =    231
1bar, and number of water molecules = 9, 548 for Cys and 9,549 for Met systems.    232

## 3.3    Tripeptide Subsystem    233

### 3.3.1    Description of Model System    234

The tripeptide active site system (His–Asn–Cys) is highlighted in the view of the    235
full intein crystal structure in Fig. 3b. Gas phase calculations were used to study    236
the effect of site-directed mutagenesis (see Fig. 4). Intein crystal structures usually    237
include a hydrogen bond between the $N^{\delta}$–H of the (penultimate) His side chain and    238
the carbonyl O of Asn, the final amino acid of the intein [45–49]. Although the    239
penultimate intein His residue has been previously assumed to be the proton donor    240
for C-terminal cleavage reaction in the context of splicing [50], further inspection    241
revealed that this was not the case for pH dependent C-terminal cleavage. For a    242
simple proton-catalyzed reaction, there is an inverse linear rate dependence on the    243
pH, which was observed experimentally for the C-terminal cleavage reaction [42].    244
Since the ability of His to act as an acid is based on its local $pK_a$ value, the expected    245
pH-rate curve should be non-linear, specifically sigmoidal in shape, which is in    246
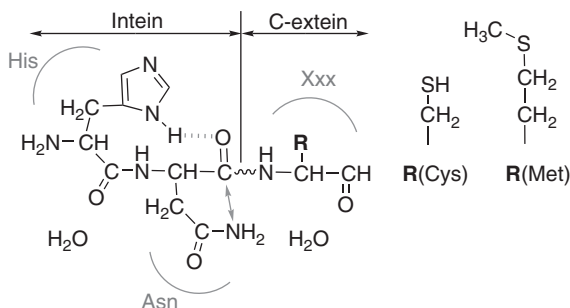contrast to the linearity observed experimentally.    247

**Fig. 4** The C1-block His–Asn–Xxx active site is shown. The highly conserved H-bond is shown with a dotted line, the cyclization coordinate of Asn is shown with an arrow, and the scissile peptide bond is shown with a wavy line. Side chains for Cys and Met are shown, although Ala, Val, Thr and Ser were also considered

The proposed N-protonation mechanism begins with the protonation of the peptide N by a hydronium ion ($H_3O^+$). This in turn causes the scissile peptide bond to elongate, and hence reduces the energy necessary for peptide bond cleavage after Asn cyclization. After Asn cyclization and aminosuccinimide formation, the extra proton passes to the cleaved C-extein N-terminus ($-NH_2$), which is excised and leaves with a positive charge ($-NH_3^+$. Although O-protonation was more energetically favorable for a generic or average peptide that was fully solvent exposed, in the case of the intein C-terminal active site, the carbonyl O was strongly hydrogen bonded to the $N^\delta$–H of His and was also pointed inward, toward the core of the protein and away from the main body of solvent. The Asn cyclization reaction after O-protonation instead of N-protonation has been shown to require more energy and does not lead to cleavage of the peptide bond [12].

Prior to the QM/MM full protein study, the His–Asn–Cys tripeptide system (Fig. 4) was studied with an isolated gas phase reaction.[2] Certain constraints were included to ensure that the backbone structure reflects that of the protein crystal structure: both terminal backbone atoms were geometrically fixed in the crystal structure configuration, both dihedral angles are constrained to values from the crystal structure and throughout the classical molecular dynamic trajectories, and the hydrogen bond between $N^\delta$–H of His and the carbonyl O of Asn was constrained at a distance of 1.8 Å. Without these constraints, the subsystem would likely rearrange into a structure that does not represent the intein C-terminal structure but does minimize the gas phase energy. By contrasting the effects of mutations, electronic structure properties at critical points were studied, including those at the purely quantum mechanical transition state.

---

[2] Gas phase energy barriers are typically higher than barriers that include electrostatic contributions such as implicit solvent calculations.

### 3.3.2   Energetic Results

For the N-protonation mechanism calculated with the tripeptide system, the com- 273
putational energy barrier for the His–Asn–Cys system in the gas phase was 27.95 274
kcal/mol, in good agreement with the experimental results of ∼21 kcal/mol [42]. For 275
a system roughly 30 atoms smaller, the previous gas phase energy barrier was ∼33 276
kcal/mol [12]. This difference indicates that even the most basic approximation of 277
the tertiary structure is important for accurate prediction of certain reaction energy 278
barriers, as we will see with the QM/MM reaction. Additionally, we have tested and 279
confirmed that the hydrogen bond between $N^\delta$–H of His and the carbonyl O of Asn 280
(dashed line in Fig. 4) caused O to not accept a proton from $H_3O^+$. This hydrogen 281
bond is usually found at the C-terminus of inteins and is important for reducing the 282
possibility of proton transfer to the carbonyl O. In fact, the normally highly exother- 283
mic reaction for $H_3O^+$ to donate a proton to the carbonyl O atom is endothermic 284
for cases where O is hydrogen bonded with another group [51]. 285

Table 1 summarizes the calculated energy barriers and relative rate constants 286
for the gas phase tripeptide system with several His–Asn–Xxx mutations. By in- 287
cluding additional atoms, the gas phase energy barrier with Xxx = Cys (27.95 288
kcal/mol) was less than the previously calculated barrier for a smaller system (33 289
kcal/mol [12]) due to polarity and geometrical effects. The larger system used here 290
was expected to more closely match the experiment of 21 kcal/mol, which it does, 291
because of the additional mechanical and electronic influences of nearby protein 292
and solvent groups. 293

The energy barrier of the His–Asn–Met system was 1.63 kcal/mol higher than the 294
His–Asn–Cys system, which corresponds to a 5.83% increase in the energy barrier. 295
When Cys was mutated to Met, the relative C-terminal reaction rate was predicted 296
to be 0.07 as fast, or decreased by more than an order of magnitude (14.0), which 297
is consistent with experimental results [42, 43]. Interestingly, this model predicts 298
that Thr and Ser instead of Cys will be slightly more effective at pH-dependent 299
C-terminal cleavage, a prediction that is consistent with the +1 position being oc- 300

t1.1 **Table 1**  Tripeptide energy barriers (ΔE) for various C-extein mutations (His–Asn–Xxx), percent change (%ΔE) from His–Asn–Cys energy barrier, and expected change in reaction rate $k_{rel}$ compared to His–Asn–Cys. Structures were geometrically optimized with the B3LYP/6-311++G(d,p) level of theory. The percent change in the energy barrier, $\%\Delta E \equiv \frac{\Delta E_{Xxx} - \Delta E_{Cys}}{\Delta E_{Xxx}} * 100\%$. Reaction rates $k$ are relative to the His–Asn–Cys wildtype at $T = 310.15$ K (37 °C). The Arrhenius equation was used to compare the relative reaction rates between two mutants: $k = k_1/k_2 = e^{-(\Delta E_1 - \Delta E_2)/RT}$, where $k_i$ and $\Delta E_i$ were the reaction rate and energy barrier for the $i^{th}$ mutant, respectively; $R$ was the gas constant and $T$ was the temperature in Kelvin

| t1.2 | Mutant (Xxx) | ΔE (kcal/mol) | %ΔE | $k_{rel}$ |
|---|---|---|---|---|
| t1.3 | **Cys** | **27.95** | 0.00 | **1** |
| t1.4 | Thr | 27.56 | −1.39 | 1.88 |
| t1.5 | Ser | 27.75 | −0.71 | 1.38 |
| t1.6 | Ala | 28.64 | 2.46 | 0.32 |
| t1.7 | Val | 28.97 | 3.64 | 0.19 |
| t1.8 | **Met** | **29.58** | **5.83** | **0.07** |

P.T. Shemella and S.K. Nayak

cupied by Cys, Thr, or Ser in nature, and will be tested in experiment. In the context 301
of splicing, experiments have shown that Cys, Ser, and Thr are the only amino acids 302
with the ability to complete the transesterification step of splicing [5], which is con- 303
sistent because they also are the most efficient at C-terminal cleavage according to 304
the calculations presented here. 305

### 3.3.3 Charge Analysis 306

Natural Populations Analysis (NPA) [52] was used to study the electron population 307
and the partial atomic charges. Figure 5a illustrates the effect of amino acid mutation 308
on the scissile peptide bond distance and Fig. 5b shows the sum of the NPA charges 309
for the mutated C-extein residue, starting with the -NH at the scissile junction and 310
including the side chain. The scissile bond distance and charge results are shown 311
as a function of each mutant's energy barrier, and include the normal amide, the 312
N-protonated amide, and the transition state corresponding to the pH dependent C- 313
terminal cleavage reaction. For the neutral amide, the C–N scissile peptide bond 314
distance was 1.3492 Å for Cys, which decreased to 1.3455 Å for Met. Although this 315
change was extremely small, it does confirm that the amino acid side chain played a 316
small but perceptible role in the properties of a normal peptide bond (which is well 317
known from proton exchange experiments [53]). For the N-protonation step and 318
then the Asn cyclization transition state, the correlation between short scissile bond 319
distance and high energy barrier was more apparent: a shorter peptide bond implied 320
more $\pi$-bond resonance between C and N, less $\pi$-bond resonance between C and O, 321
and more energy was required to break the C–N bond. An elongated peptide bond 322
implied less $\pi$ bonding between C and N and less energy necessary for peptide bond 323
cleavage [54]. 324

A correlation between the energy barrier and the net charge can be seen (Fig. 5b), 325
especially for the Cys/Met mutation, signifying that the residues that were able to 326
accept more electrons exhibit a reduced energy barrier whereas the residues that 327
were less likely or unable to accept electrons displayed an increased energy barrier. 328

## 3.4 Single Amino Acid Molecules 329

### 3.4.1 Electron Affinity and Ionization Potential Analysis 330

To further elucidate the effect of the mutation of the first C-extein amino acid side 331
chain on the energy barrier, the isolated Cys and Met amino acids were studied. The 332
electron affinities (EA) and ionization potentials (IP) for each were calculated with 333
the B3LYP/6-311++G(d,p) level of theory. The EA for Cys, (the amount of energy 334
gained or lost when the system goes from neutral to negatively charged), was 6.79 335
kcal/mol. For Met, the EA was 8.27 kcal/mol, signifying that the side chain of the 336
gas phase Cys residue was more electronegative than for Met. The reason that Cys 337
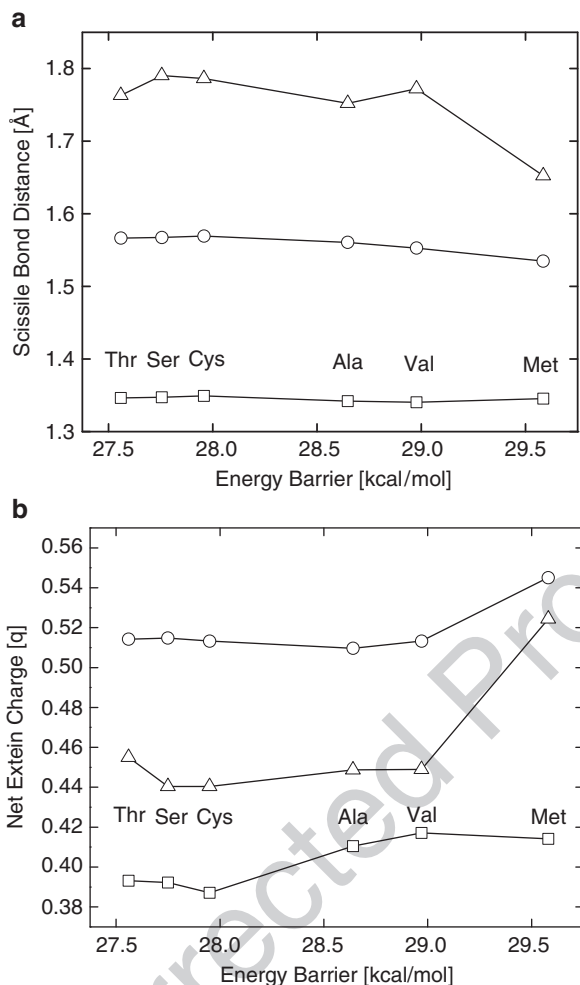
**Fig. 5** Relaxed scissile peptide bond distance (**a**) and NPA charges summed for atoms on the C-extein (**b**) for the tripeptide gas phase system, His–Asn–Xxx (Xxx = Thr, Ser, Cys, Ala, Val, Met). Both the scissile bond distance and the net charge for the C-extein amino acid (Xxx) are plotted as a function of the specific mutant's energy barrier and are shown for the normal amide, (□); the N-protonated amide, (○); and the Asn cyclization transition state (△)

was more stable with charge than Met was due to the bonding for each S atom. Al-  338
though each side chain contained an S atom, for Cys the S atom was bonded to one  339
methyl group and one H atom. For Met, both bonds of the S atom were to methyl  340
groups, hence different electron occupation properties. In changing from neutral  341
to negatively charged, the partial charge of S for Cys changed from −0.01051 to  342
−0.11874 units of charge, corresponding to the addition of 0.10823 electrons. For  343
Met, the charge went from 0.16894 to 0.12532 units of charge, corresponding to  344

the gain of only 0.04362 electrons. The S of Cys was able to accommodate more    345
than twice the amount of delocalized electron population as compared to Met, indi-    346
cating more energetic stability in the negatively charged system. The difference in    347
ionization potential (IP) for the same isolated Cys and Met amino acids was calcu-    348
lated. The removal of one electron from Cys required 203.05 kcal/mol while that    349
for Met was 191.14 kcal/mol. Combining the fact that Met was more stable when an    350
electron was removed, and the fact that Cys was more stable when an electron was    351
added, we conclude that the "electron pulling" and "electron pushing" properties of    352
the first C-extein amino acid side chain must have an effect on the actual properties    353
of the scissile peptide bond.    354

### 3.4.2  Energetic Analysis of Molecular Orbitals near the Fermi Energy    355

For the isolated amino acids (Thr, Ser, Cys, Ala, Val, and Met), the highest occupied    356
molecular orbital (HOMO) for the neutrally charged system as well as the negatively    357
charged system was compared. The difference in energy between the HOMO of the    358
electron doped (negatively charged) and the neutral system is termed the energy gap,    359
and is shown in Fig. 6. From this analysis of the negatively charged amino acids (ge-    360
ometrically optimized with neutral charge), the isolated amino acids are ranked in    361
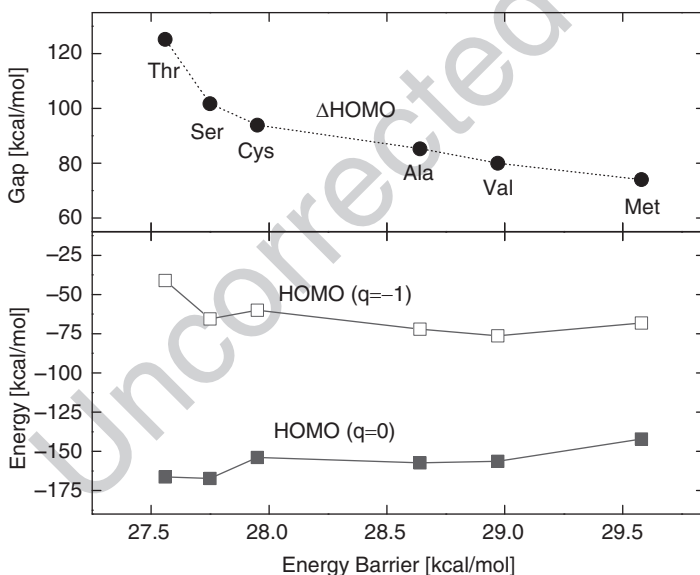


**Fig. 6** Energies for the highest occupied molecular orbital (HOMO) for the neutral system (■) and the negatively charged system (□) for the isolated amino acid molecules (Thr, Ser, Cys, Ala, Val, Met), shown in order of their energy barrier found independently for the tripeptide reaction calculation. The difference between these energies is the energy gap (●) and is clearly dependent to the energy barrier for the given mutant

order of the energy barrier found when they are the mutant for the tripeptide system,    362
and there was a clear trend in the energy gap between the neutral and negatively    363
charged molecules. The energy gap was closely related to the electron affinity of the    364
molecule: as the energy barrier increased for a particular mutant, the gap decreased.    365
This single amino acid analysis is of particular interest because from the electronic    366
structure properties of an isolated molecule representing an amino acid side chain,    367
calculated properties such as the electron affinity, the ionization potential, and the    368
molecular orbital energy levels may explain and perhaps predict the relative reaction    369
rate for an unknown mutant at the first C-extein position.    370

The localization of the EA densities found for molecules characterized in Fig. 6    371
is plotted as a volumetric surface in Fig. 7, which shows the difference in electron    372
density between the neutral (optimized geometry) and negatively charged (single    373
point geometry) single amino acid residues (Thr, Ser, Cys, Ala, Val, and Met). The    374
presence of electrons on the molecular side chain was observed for amino acids that    375
are more efficient when downstream of the scissile peptide bond in intein C-terminal    376
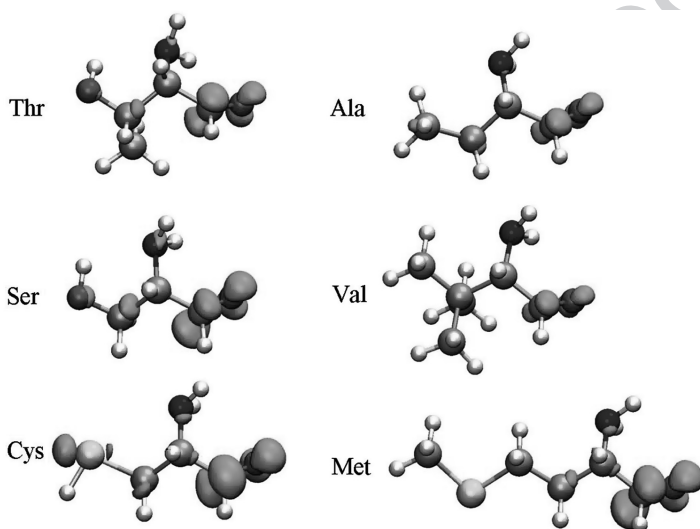cleavage.    377



**Fig. 7** The electron affinity (EA) density for single amino acid molecules (Thr, Ser, Cys, Ala, Val, and Met). The electron density surface describes the delocalization of the electron affinity when an electron is added to the system, thus going from neutral to negatively charged ($\Delta\rho$). For downstream amino acids that were efficient at C-terminal cleavage (Thr, Ser, Cys), the EA density extended to the side chain. For amino acids that were less efficient (Ala, Val, Met), the EA density remained on the peptide-like part of the molecule, and away from the side chain. Atom colors are as follows: carbon is cyan, nitrogen is blue, oxygen is white, sulfur is yellow, and hydrogen is white; the electron density surface is green [55]

### 3.4.3  Tripeptide Analysis                                                378

Returning to the tripeptide system shown in Fig. 4, Table 2 shows electron popu-   379
lation analysis for orbitals with $l = 1$ angular momentum (2s orbital), as well as   380
total occupation for $l = 0, 1$ (2s and 2p orbitals). From the analysis of target atoms   381
belonging to the scissile peptide bond, the expected differences in electron popula-   382
tion between Cys/Met mutants were observed. Specifically, the N atom for Met was   383
generally more occupied with electrons than Cys, which gave it a greater negative   384
charge.                                                                        385

  For both mutants, the N atom showed a considerable increase of 2s electrons,   386
which corresponded to C and other atoms returning $\sigma$ electrons to N when the C-N   387
bond was elongated after N-protonation. A similar situation with $\sigma$ electron back-   388
transfer to N was found for peptide bond rotation, where at the transition state of   389
90° the N atom lost $\pi$ electrons although there was an increase in $\sigma$ electrons to N   390
[54]; this phenomenon explains why N actually became more negative as similarly   391
seen in the present study. The 2p orbitals for N showed distinct differences for the   392
Cys/Met mutations – even for the neutral ground state which was a normal amide   393
system, a distinction that signified the side chains of adjacent amino acids were   394
important in dictating the exact properties of the peptide bond.                395

  For the normal amide, the charge of the peptide N for Cys was −0.616 and for   396
Met the charge was −0.641. For the N-protonation case, the charge of N for the Cys   397
case was −0.660, where for Met the charge was −0.710. For the transition state, the   398
charge on N for Cys was −0.684, and for Met was −0.699. For all three cases the   399
charge of N for Met was more negative than for Cys, which was consistent with the   400
electron affinity calculation described previously. The side chain plays a subtle yet   401
important role in the electrostatic environment during the cleavage reaction. By hav-   402
ing less charge on N, the -NH$_2$ group is more energetically favored to leave. From   403
this electron population analysis, differences in the electronic structure of the scis-   404

t2.1  **Table 2**  Atomic orbital populations for the 2s and net 2p orbitals as well as the total electronic occupation for the peptide N atom in the gas phase tripeptide calculation. N is generally less occupied by electrons for Cys as compared to Met, which is consistent with single amino acid electron affinity results. The sum of electron occupation for the $2p_x$, $2p_y$, and $2p_z$ orbitals is written as 2p. The NPA charge is calculated by subtracting the total electron occupation from the atomic number; a larger electron occupation signifies a more negative charge

| t2.2 t2.3 | | | Occupation | | |
|---|---|---|---|---|---|
| t2.4 | Orbital | Mutant | Neutral ground state | N-protonated | Transition state |
| t2.5 | [2s] | Cys | 1.250 | 1.359 | 1.386 |
| t2.6 | | Met | 1.259 | 1.360 | 1.376 |
| t2.7 | [2p] | Cys | 4.341 | 4.285 | 4.277 |
| t2.8 | | Met | 4.357 | 4.329 | 4.299 |
| t2.9 | Total | Cys | 7.616 | 7.660 | 7.684 |
| t2.10 | | Met | 7.641 | 7.710 | 7.699 |

sile peptide bond for Cys and Met were observed, which explained why the energy   405
barrier for Cys and Met mutants would be distinct despite an identical mechanism.   406

# 4  Reaction Analysis with QM/MM Calculations                                     407

The full protein QM/MM reaction profile was initially calculated with the QM ac-   408
tive site region of His–Asn–Cys, and two water molecules (2346 protein atoms,   409
4161 water atoms, and total 53 QM atoms) [56]. Figure 8 shows the QM/MM   410
energy barrier with and without electrostatic embedding. The energy barrier was   411
24.96 kcal/mol for the QM/MM calculation with geometry optimization, in excel-   412
lent agreement with the 21 kcal/mol measured experimentally [42].                413

## 4.1  *Effect of Mutation on Energy Barriers*                                    414

The energy barrier difference for the Cys/Met mutation is of interest in the context   415
of a QM/MM calculation, but because the Met side chain was too spatially extended   416
to simply replace the smaller Cys side chain, additional classical MD simulations   417
were performed (starting from the initial intein plus extein structure) but with Met   418
at the C-extein +1 residue. Once the full protein system was equilibrated, the QM   419
active site was partitioned to be His–Asn–Met plus the two water molecules in the   420
same location as before (59 total QM atoms). The Asn cyclization reaction coor-   421
dinate was scanned after N-protonation by $H_3O^+$. To compare the effect of the   422
Met/Cys mutation directly, the smaller Cys was substituted for Met, and the ge-   423
ometry was again relaxed. By doing this, the change in reaction energies may be   424
compared directly because the original protein structures were common for both   425
Met and Cys residues.                                                            426

   These structures were in near total overlap, with the exception of the side chain   427
of the (+1) amino acid, either -$CH_2$-SH for Cys, or -$(CH_2)_2$-S-$CH_3$ for Met. Us-   428
ing the B3LYP/6-31G(d,p) level of theory, independent reaction profiles for the   429
Met/Cys mutation were calculated. For Met the barrier was 27.07 kcal/mol and for   430
Cys was 26.17 kcal/mol. The His–Asn–Met QM active site (as part of the QM/MM   431
system) had an energy barrier of 0.90 kcal/mol higher than His–Asn–Cys, which   432
corresponded to ratio between reaction rates of $k = k_{Cys}/k_{Met} = 0.22$, in good   433
agreement with experimental results and consistent with the tripeptide system con-   434
clusions [42, 43].                                                               435
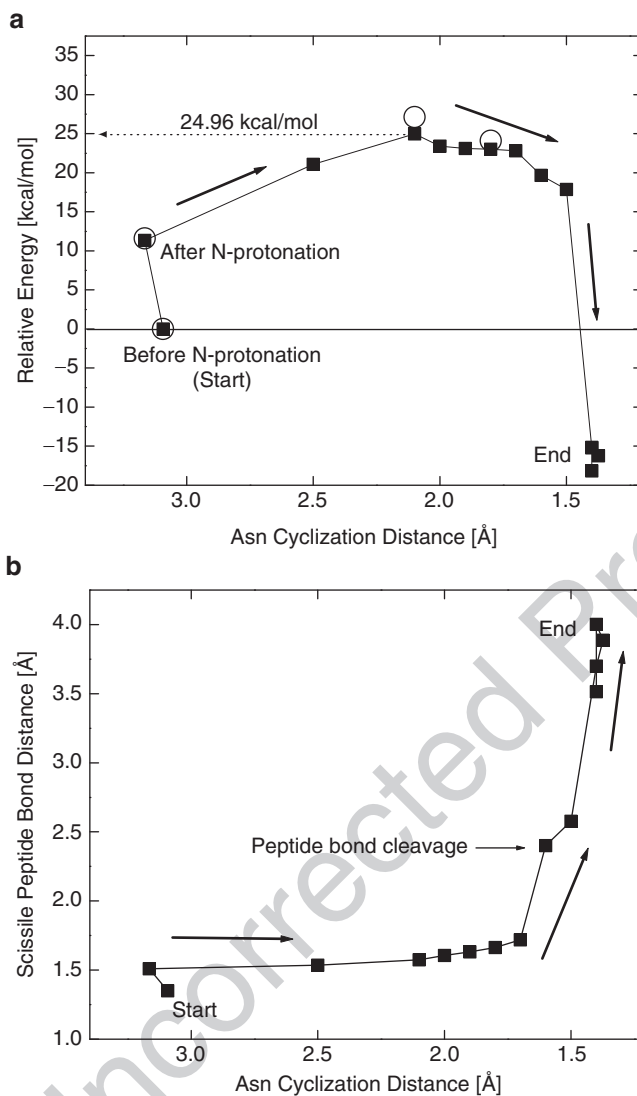
**Fig. 8** Combined QM/MM reaction energy profile (**a**) and distance of the scissile peptide bond during breakage (**b**) for His–Asn–Cys plus two water QM system. QM/MM geometry optimization (■). QM/MM + charge embedding single point energies (○)

## *4.2   Effect of Mutation on Electron Occupation*                                    436

In addition to energy barriers, the Mulliken charge [57] was calculated for critical   437
atoms.[3] For the N atom of the scissile bond and for the ground state, the partial     438
charge was $-0.538$ for Cys and for Met was $-0.545$. For the N-protonation state the   439
partial charge of N was $-0.609$ for Cys and was $-0.615$ for Met. At the transition    440
state, the charge for Cys was $-0.584$ and for Met was $-0.598$. In all cases the       441
partial charge of the N atom for the Met mutant was more negative, which was            442
consistent with the tripeptide results, and is explained by using the electron affinity 443
and ionization potential for the isolated Cys and Met amino acids. When the net         444
Mulliken charge was summed for the C-extein residue (Cys or Met) in the QM/MM           445
context for the normal amide ground state, for Met the net charge was 0.225, and        446
for Cys the net charge was 0.209.                                                       447

Within the QM/MM system, the charge for the backbone and side chain of the              448
first C-extein residue was added. The net charge of Cys was more negative than Met,     449
which is in agreement with the model QM calculations described in the preceding         450
paragraphs.                                                                             451

By combining model system QM calculations and full-protein QM/MM simula-               452
tions, the non-mechanistic regulation of reaction rate regulation for single amino      453
acid mutations near to the active site was confirmed, explained, and predicted.         454
Similar methods are also useful for testing an unknown mechanism based on the           455
correlated experimental results of kinetic data (from non-essential amino acid site-    456
directed mutagenesis).                                                                  457

## 5   Conclusions                                                                      458

The C-terminal cleavage reaction and the previously proposed N-protonation mech-        459
anism were tested by increasing the QM system size by 30 atoms to at least 53           460
atoms. In addition, full-protein QM/MM analysis was performed. The pH dependent         461
C-terminal cleavage reaction undergoes simple proton-catalysis by a hydronium ion       462
that protonates the peptide N atom. The peptide bond, now resonance destabilized,       463
is elongated and the peptide C atom is open for attack by the Asn side chain. Dur-      464
ing Asn cyclization, the peptide bond cleaves while an aminosuccinimide ring is         465
formed. The final step involves the donation of the extra proton on the aminosuccin-    466
imide to the -NH$_2$ leaving group *via* water, thus making the leaving group positively 467
charged. Our QM/MM results included the effects from the protein interior, both         468
mechanical and electrostatic.                                                           469

The "non-mechanistic" role of the first amino acid of the C-extein was confirmed.       470
This amino acid, although not necessary for C-terminal cleavage, did have an effect     471
on the reaction rate by about an order of magnitude, as measured by Wood et al.         472

---

[3] Natural Population Analysis (NPA) is not implemented with QM/MM at this time.

[42, 43, 58]. In this study, the precise energy barrier for C-terminal cleavage (and hence reaction rate) was shown to be dependent on the side chain of the amino acid downstream from the scissile bond. Explained by the electron occupation and partial atomic charges for each residue at the C +1 position, considerable differences that led to a distinction in energy barriers were calculated and found to be in agreement with experimentally observed reaction rates.

# References

1. Belfort, M., Stoddard, B., Wood, D., Derbyshire, V.: Homing Endonucleases and Inteins (Nucleic Acids and Molecular Biology). Springer, Heidelberg (2005)
2. Perler, F., Davis, E., Dean, G., Gimble, F., Jack, W., Neff, N., Noren, C., Thorner, J., Belfort, M.: Nucl. Acid Res. **22**(7), 1125 (1994)
3. Hiraga, K., Derbyshire, V., Dansereau, J., Van Roey, P., Belfort, M.: J. Mol. Biol. **354**(4), 916 (2005)
4. Shingledecker, K., Jiang, S.Q., Paulus, H.: Arch. Biochem. Biophys. **375**(1), 138 (2000)
5. Paulus, H.: Ann. Rev. Biochem. **69**(1), 447 (2000)
6. Miao, J., Wu, W., Spielmann, T., Belfort, M., Derbyshire, V., Belfort, G.: Lab Chip **5**(3), 248 (2005)
7. Banki, M., Feng, L., Wood, D.: Nat. Methods **2**(9), 659 (2005)
8. Paulus, H.: Front. Biosci. **8**, s1157 (2003)
9. Mootz, H., Muir, T.: J. Am. Chem. Soc. **124**(31), 9044 (2002)
10. Muralidharan, V., Muir, T.: Nat. Methods **3**, 429 (2006)
11. Senn, H., Thiel, W.: Angew. Chem. Int. Edit. **48**(7) (2009)
12. Shemella, P., Pereira, B., Zhang, Y., Van Roey, P., Belfort, G., Garde, S., Nayak, S.: Biophys. J. **92**(3), 847 (2007)
13. Shemella, P., Pereira, B., Van Roey, P., Belfort, G., Garde, S., Nayak, S.: Controlling C-terminal cleavage rate of an intein through extein mutation: A quantum mechanical study. (submitted)
14. Hohenberg, P., Kohn, W.: Phys. Rev. B **136**(3B), 864 (1964)
15. Kohn, W., Sham, L.: Phys. Rev. **140**(4A), 1133 (1965)
16. Becke, A.: J. Chem. Phys. **98**, 5648 (1993)
17. Becke, A.: Phys. Rev. A **38**(6), 3098 (1988)
18. Lee, C., Yang, W., Parr, R.: Phys. Rev. B **37**(2), 785 (1988)
19. Perdew, J., Wang, Y.: Phys. Rev. B **45**(23), 13244 (1992)
20. Frisch, M.J., Trucks, G.W., Schlegel, H.B., Scuseria, G.E., Robb, M.A., Cheeseman, J.R., Montgomery, Jr. J.A., Vreven, T., Kudin, K.N., Burant, J.C., Millam, J.M., Iyengar, S.S., Tomasi, J., Barone, V., Mennucci, B., Cossi, M., Scalmani, G., Rega, N., Petersson, G.A., Nakatsuji, H., Hada, M., Ehara, M., Toyota, K., Fukuda, R., Hasegawa, J., Ishida, M., Nakajima, T., Honda, Y., Kitao, O., Nakai, H., Klene, M., Li, X., Knox, J.E., Hratchian, H.P., Cross, J.B., Bakken, V., Adamo, C., Jaramillo, J., Gomperts, R., Stratmann, R.E., Yazyev, O., Austin, A.J., Cammi, R., Pomelli, C., Ochterski, J.W., Ayala, P.Y., Morokuma, K., Voth, G.A., Salvador, P., Dannenberg, J.J., Zakrzewski, V.G., Dapprich, S., Daniels, A.D., Strain, M.C., Farkas, O., Malick, D.K., Rabuck, A.D., Raghavachari, K., Foresman, J.B., Ortiz, J.V., Cui, Q., Baboul, A.G., Clifford, S., Cioslowski, J., Stefanov, B.B., Liu, G., Liashenko, A., Piskorz, P.,

AQ3

Komaromi, I., Martin, R.L., Fox, D.J., Keith, T., Al-Laham, M.A., Peng, C.Y., Nanayakkara, A., Challacombe, M., Gill, P.M.W., Johnson, B., Chen, W., Wong, M.W., Gonzalez, C., Pople, J.A.: Gaussian 03, Revision C.02. Gaussian, Inc., Wallingford, CT (2004)

21. Zhang, X., Zhang, X., Bruice, T.: Biochemistry **44**(31), 10443 (2005)

22. Zhang, X., Bruice, T.: Proc. Natl. Acad. Sci. **103**(44), 16141 (2006)

23. Moeller, C., Plesset, M.: Phys. Rev. **46**(7), 618 (1934)

24. Head-Gordon, M., Pople, J., Frisch, M.: Chem. Phys. Lett. **153**(6), 503 (1988)

25. Frisch, M., Head-Gordon, M., Pople, J.: Chem. Phys. Lett. **166**(3), 275 (1990)

26. Frisch, M., Head-Gordon, M., Pople, J.: Chem. Phys. Lett. **166**(3), 281 (1990)

27. Miehlich, B., Savin, A., Stoll, H., Preuss, H.: Chem. Phys. Lett. **157**(3), 200 (1989)

28. Hehre, W., Ditchfield, R., Pople, J.: J. Chem. Phys. **56**(5), 2257 (1972)

29. McLean, A., Chandler, G.: J. Chem. Phys. **72**(10), 5639 (2006)

30. Miertus, S., Scrocco, E., Tomasi, J.: Chem. Phys. **55**(11), 117 (1981)

31. Cances, M., Mennucci, B., Tomasi, J.: J. Chem. Phys. **107**(8), 3032 (1997)

32. Van Roey, P., Pereira, B., Li, Z., Hiraga, K., Belfort, M., Derbyshire, V.: J. Mol. Biol. **367**(1), 162 (2007)

33. Davis, E., Sedgwick, S., Colston, M.: J. Bacteriol. **173**(18), 5653 (1991)

34. Cornell, W., Cieplak, P., Bayly, C., Gould, I., Merz, K., Ferguson, D., Spellmeyer, D., Fox, T., Caldwell, J., Kollman, P.: J. Am. Chem. Soc. **117**(19), 5179 (1995)

35. Van der Spoel, D., Lindahl, E., Hess, B., Groenhof, G., Mark, A., Berendsen, H.: J. Comput. Chem. **26**(16), 1701 (2005)

36. Maseras, F., Morokuma, K.: J. Comput. Chem. **16**(9), 1170 (1995)

37. Vreven, T., Morokuma, K., Farkas, O., Schlegel, H., Frisch, M.: J. Comput. Chem. **24**(6), 760 (2003)

38. Gao, J., Truhlar, D.: Ann. Rev. Phys. Chem. **53**, 467 (2002)

39. Cui, Q., Elstner, M., Karplus, M.: J. Phys. Chem. B **106**(10), 2721 (2002)

40. Torrent, M., Vreven, T., Musaev, D., Morokuma, K., Farkas, O., Schlegel, H.: J. Am. Chem. Soc. **124**(2), 192 (2002)

41. Dellago, C., Bolhuis, P., Csajka, F., Chandler, D.: J. Chem. Phys. **108**(5), 1964 (1998)

42. Wood, D., Derbyshire, V., Wu, W., Chartrain, M., Belfort, M., Belfort, G.: Biotechnol. Progr. **16**(6), 1055 (2000)

43. Wood, D.: Generation and application of a self-cleaving protein linker for use in single-step affinity fusion based protein purification. Ph.D. thesis, Rensselaer Polytechnic Institute (2000)

44. Wu, W., Wood, D., Belfort, G., Derbyshire, V., Belfort, M.: Nucl. Acid Res. **30**(22), 4864 (2002)

45. Mizutani, R., Nogami, S., Kawasaki, M., Ohya, Y., Anraku, Y., Satow, Y.: J. Mol. Biol. **316**(4), 919 (2002)

46. Poland, B., Xu, M., Quiocho, F.: J. Biol. Chem. **275**(22), 16408 (2000)

47. Klabunde, T., Sharma, S., Telenti, A., Jacobs, W., Sacchettini, J.: Nat. Struct. Biol. **5**(1), 31 (1998)

48. Duan, X., Gimble, F., Quiocho, F.: Cell **89**(4), 555 (1997)

49. Ichiyanagi, K., Ishino, Y., Ariyoshi, M., Komori, K., Morikawa, K.: J. Mol. Biol. **300**(4), 889 (2000)

50. Ding, Y., Xu, M., Ghosh, I., Chen, X., Ferrandon, S., Lesage, G., Rao, Z.: J. Biol. Chem. **278**(40), 39133 (2003)

51. Shemella, P., Nayak, S.: (unpublished)

52. Reed, A., Curtiss, L., Weinhold, F.: Chem. Rev. **88**(6), 899 (1988)

53. Bai, Y., Milne, J., Mayne, L., Englander, S.: Prot. Struct. Funct. Genet. **17**(1), 75 (1993)

54. Milner-White, E.: Prot. Sci. **6**(11), 2477 (1997)

14. Humphrey, W., Dalke, A., Schulten, K.: J. Mol. Graph. **14**(1), 33 (1996)

56. Pereira, B., Jain, S., Garde, S.: J. Chem. Phys. **124**, 074704 (2006)

57. Mulliken, R.: J. Chem. Phys. **23**(10), 1833 (1955)

58. Wood, D., Wu, W., Belfort, G., Derbyshire, V., Belfort, M.: Nat. Biotechnol. **17**(9), 889 (1999)

AQ4

# AUTHOR QUERIES

AQ1.  Please provide keywords.
AQ2.  Please sepcify the corresponding author.
AQ3.  Please update Shemella et al.
AQ4.  Please provide year for Ref. Shemella and Nayak.